

Do mathematicians dream of data?

Katja Berčič
University of Ljubljana and IMFM

LMFDB, Computation, and Number Theory (LuCaNT) 2025
ICERM, July 7-11

No.*

* The rest of this talk is about the asterisk

No: most (but not all) mathematicians (perhaps those that don't call themselves applied) would say so. There are some subtleties to consider, and this talk is about those.

1. Overview
A trilogy in four parts

2. What is data (in math)?
Terminology, history and the present

3. What can we do better?
Data management, a TL;DR

4. Where is it & what is it like?
Towards a survey

5. How can we trust it?
Validation and correctness



An incomplete
list of my influences

On the last slide, you'll get a QR code to a page with a list of references (papers and otherwise), and if I had a productive evening, the slides, too.

Disclaimer

Before we get to the good part, I have a disclaimer to make. The conversation around data in mathematics is still in its early days, and so are the words we use to talk about it. My perspective is inevitably shaped by my background in combinatorics and graph theory — though it seems to me that everyone else is grappling with similar predicaments.

<p>1. Overview</p> <p>A trilogy in four parts</p> <p>2. What is data (in math)?</p> <p><i>Terminology, history and the present</i></p> <p>3. What can we do better?</p> <p><i>Data management: a TL;DR</i></p> <p>4. Where is it & what is it like?</p> <p><i>Towards a survey</i></p> <p>5. How can we trust it?</p> <p><i>Validation and correctness</i></p>	
---	--

A long time ago, when I was a young and naive undergraduate student, I built a database of tournaments for the Slovenian Go Association. Soon after that I started my PhD and one of the first things I encountered was my advisor's dataset of cubic vertex-transitive graphs. To me, it seemed rather urgent that something like that belonged in a database, not just as a plain text file.

Here's how this talk will unfold:

- As a warm-up (not one of the main sections), we'll spend a bit of time on terminology — with some neat pictures of historical data.
- Then we'll move on to the three main parts: first, FAIR&RDM (the boring bit); second, where some of the datasets are and what they look like; and finally, how much we can trust the data.
- If we're counting generously, this slide could be the “fifth part,” making this talk a trilogy in five parts. Don't panic; there won't be a sixth.

2. What is data?

A poll

Scan the QR code on the right
or go to menti.com and enter
the code 2106 7957

When answering, consider

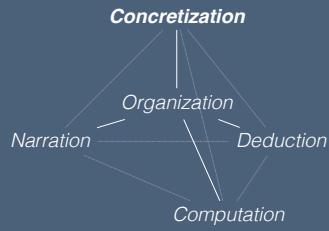
- who produces the data,
- who are the users,
- what is the content.



Language shapes thinking

A story of a word

- Which word might fit the activity of coming up with (counter)examples?
- What should be listed in Mathbases?
- Examples are important! As are collections of examples.
- How about "do I need a research data management plan for this?"



Jacques Carette, Bill Farmer, Michael Kohlhase, and Florian Rabe were searching for a word to describe the core mathematical activity involved in constructing databases. When I joined the conversation, it was still being called tabulation, but even before that, they thought that there might be a more fundamental process beneath that. We eventually settled on concretization (making abstract ideas concrete), though I'm not sure we were ever all completely satisfied with it.

At the time, I was thinking about — and worse, trying to implement — a generic math database (I like to think I came out of that a little wiser). Wrestling with that elusive fifth activity helped me organize my thoughts and I now see instances of concretization in more places than before.

Beyond listing examples

Let us remember that the real world is wide

- "Research data are all digital and analog objects generated or handled in the process of doing research" (*MaRDI*)
- *Right:* the listing on the Cornell Mathematics Library page on "Math Databases".

- MathSciNet
- Zentralblatt für Mathematik
- Google Scholar
- Wikipedia
- MacTutor History of Mathematics
- Scopus
- The Web of Science
- Mathworld
- Jahrbuch-Project Electronic Research Archive for Mathematics (mathematics literature 1868 - 1943)
- arXiv
- ERIC (index in the field of Education, including Education in Mathematics)
- Wolfram|Alpha ("this search engine allows you to enter a query and returns an answer from structured data")

"All digital and analog objects" includes: paper publications, proofs, computational results (and more). Does this mean that all mathematicians should have a research data management plan when they start writing a paper? Probably not, but perhaps they should.

Cuneiform and the Greeks

Data probably goes as far back as math

- Earliest datable table containing mathematical computations: length measurements, corresponding areas (Sumer, c. 2600 BCE).
- *picture*
A list of Pythagorean triples (Babylon, c. 1800 BCE).
- The list of the Platonic solids (Theaetetus of Athens, c. 417 - c. 369 BCE).



Bill Casselman

Eleanor Robson, "Tables and tabular formatting in Sumer, Babylonia, and Assyria, 2500 BCE-50," Campbell-Kelly et al [eds]. *The History of Mathematical Tables from Sumer to Spreadsheets* [2003]

Results of computations

Persistence: data outlasting the computation before computers

- picture

*Trigonometric, logarithmic, and exponential functions become subjects of tables
(Napier's Mirifici logarithmorum, trigonometric and log trig data for 34 degrees)*

- Math Tables Project (1938 - 1946): human computers constructed tables of mathematical functions (needed for hand computation).

Gr.	34	+	-	
1	1.0000000			1.0000000
2	0.9999999			0.9999999
3	0.9999998			0.9999998
4	0.9999997			0.9999997
5	0.9999996			0.9999996
6	0.9999995			0.9999995
7	0.9999994			0.9999994
8	0.9999993			0.9999993
9	0.9999992			0.9999992
10	0.9999991			0.9999991
11	0.9999990			0.9999990
12	0.9999989			0.9999989
13	0.9999988			0.9999988
14	0.9999987			0.9999987
15	0.9999986			0.9999986
16	0.9999985			0.9999985
17	0.9999984			0.9999984
18	0.9999983			0.9999983
19	0.9999982			0.9999982
20	0.9999981			0.9999981
21	0.9999980			0.9999980
22	0.9999979			0.9999979
23	0.9999978			0.9999978
24	0.9999977			0.9999977
25	0.9999976			0.9999976
26	0.9999975			0.9999975
27	0.9999974			0.9999974
28	0.9999973			0.9999973
29	0.9999972			0.9999972
30	0.9999971			0.9999971
31	0.9999970			0.9999970
32	0.9999969			0.9999969
33	0.9999968			0.9999968
34	0.9999967			0.9999967
35	0.9999966			0.9999966
36	0.9999965			0.9999965
37	0.9999964			0.9999964
38	0.9999963			0.9999963
39	0.9999962			0.9999962
40	0.9999961			0.9999961
41	0.9999960			0.9999960
42	0.9999959			0.9999959
43	0.9999958			0.9999958
44	0.9999957			0.9999957
45	0.9999956			0.9999956
46	0.9999955			0.9999955
47	0.9999954			0.9999954
48	0.9999953			0.9999953
49	0.9999952			0.9999952
50	0.9999951			0.9999951
51	0.9999950			0.9999950
52	0.9999949			0.9999949
53	0.9999948			0.9999948
54	0.9999947			0.9999947
55	0.9999946			0.9999946
56	0.9999945			0.9999945
57	0.9999944			0.9999944
58	0.9999943			0.9999943
59	0.9999942			0.9999942
60	0.9999941			0.9999941

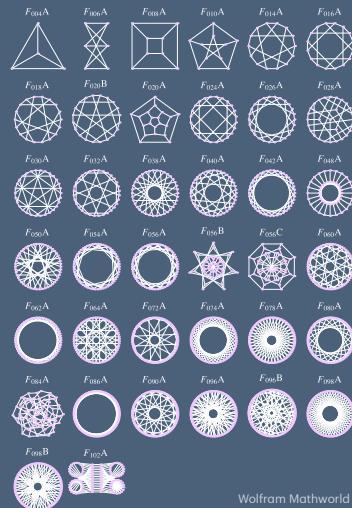
John Napier, 1614

The Foster Census

Begun in 1930, published in 1988 as a book.

Cubic symmetric graphs with up to 512 vertices.

There are 207 such graphs.
Foster only missed a handful.



Wolfram Mathworld

Symmetric graph: every ordered pair of adjacent vertices (an arc) can be mapped to any other such pair.

Modern mathematical databases
Four well-established examples, just in case

The OEIS is supported by the [On-Line Encyclopedia of Integer Sequences](#) (OEIS) Foundation.

THE ON-LINE ENCYCLOPEDIA OF INTEGER SEQUENCES

Founded in 1964 by N. J. A. Sloane

Enter a sequence, word, or sequence number:

Search [Home](#) Welcome Video

For more information about the Encyclopedia, see the [Welcome page](#).

The L-functions and modular forms database (LMFDB)

A database

The LMFDB is a database of L-functions and modular forms – primarily those associated with number fields and automorphic forms – together with the data needed to compute them.

Announcements

Coming up: LMFDB 1.0! LMFDB 1.0 is the first public release of the LMFDB. It includes many new features and improvements.

Learn more

Information available regarding the LMFDB, including how to use it and what it contains.

Documentation

Documentation for LMFDB users.

Citation and acknowledgments

How to cite the LMFDB, source code repository, acknowledgments, and references.

The House of Graphs

Search [Search History](#) Draw Graph [Meta-directory](#) Publications Help Log In Register

20 graphs found that satisfy the following conditions:

20 most recently added graphs

Download these 20 graphs in the selected format:

Graph6 [Download](#)

See also: [exact search terms](#)

HOG Id	Name	Number of Vertices	Number of Edges	Minimum Degree	Maximum Degree
53807	Asymmetric graph (332210)	7	6	5	5

π-Base

a community database of topological counterexamples

- Search spaces by name/[topo](#) or [combinations](#)
- Search spaces by properties: [non-metric](#) [continuous](#)
- Find counterexamples: [connected](#) spaces need not be path connected

Topology is a dense forest of counterexamples. A uselss map of the forest is a fine thing.
— Paraphrased from Mary Ellen Rudin's review of Counterexamples in Topology.

Questions?

3. Data management (TL;DR)

FAIR and RDMP, buzzwords of the day

- FAIR: Findable, Accessible, Interoperable and Reusable.
Not the same as open data
(*free use, accessibility, submission*)
- RDMP: Research Data Management plan



I was struggling with an image that would illustrate my frustration with some of the data in mathematics I come across (typically as some from outside of the area of mathematics), until I attempted to take a shower after checking-in on Sunday. The interface made little sense to me. There are red arrows on the right and blue arrows on the left. The main handle turns counter-clockwise. I could not get the water to be warmer than lukewarm. The smaller handle appeared to have little effect. I messed around with it for a while until, as a last resort, I turned the handle all the way into the worryingly blue zone, which unexpectedly but fortunately resulted in hot water.

<p>3. Data management (TL;DR)</p> <p>FAIR and RDMP. buzzwords of the day</p> <ul style="list-style-type: none"> • FAIR: Findable, Accessible, Interoperable and Reusable. Not the same as open data <i>(free use, accessibility, submission)</i> • RDMP: Research Data Management plan 	<p>List of the 17 representatives of IC(6,3), ordered by the RevLex-Index. The numbers above the signs indicate the elements of the corresponding basis.</p> <pre> 11121121231121231234 22332334442334445555 34445555566666666666 IC(6,3, 1) = ++++++-----+---+ IC(6,3, 2) = ++++++-----+---+ IC(6,3, 3) = ++++++-----+---+ IC(6,3, 4) = ++++++-----+---+ IC(6,3, 5) = 0+++++-----+---+ IC(6,3, 6) = 0+++++-----+---+ IC(6,3, 7) = 0+++++-----+---+ IC(6,3, 8) = 0+++++-----+0 IC(6,3, 9) = 0++++++0+++++---+ IC(6,3,10) = 0++++++0+++++---+ IC(6,3,11) = 0++++++0+++++---+ IC(6,3,12) = 0++++++0+++++0---+ IC(6,3,13) = 0++++++0+++++0+0- IC(6,3,14) = 0++++++0+++++0-+--+ IC(6,3,15) = 0000+++++-----+0+++ IC(6,3,16) = 0000+++++-----+0+++ IC(6,3,17) = 0000000000+++++---+ </pre>
--	--

The first time I saw the data on the right side of the slide, I had no idea what I was looking at. Of course, like with the tap, it is possible to figure it out (I presume), given enough time.

The FAIR guiding principles were published in 2016 and are an attempt to describe how usable data look like in very general terms.

Open in "open access" refers to the removal of financial, legal and technical barriers to data, while accessibility in FAIR refers to the data being retrievable by humans and machines.

The FAIR guidelines

Mostly about metadata: deep FAIR proposed term for corresponding properties for objects

- **Findable:** globally unique, persistent IDs, rich metadata, indexing
easy to identify and find for both humans and computers, e.g. with metadata that facilitate searching for specific datasets
- **Accessible:** stored long term, accessible and/or downloadable with well-defined access conditions, whether at the level of metadata, or at the level of the actual data.
- **Interoperable:** FAIR knowledge representation language
ready to be combined with other datasets by humans or computers, without ambiguities in the meanings of terms and values.
- **Reusable:** clear usage license, provenance, domain-relevant standards, comprehensive and relevant attributes
ready to be used for future research and to be further processed using computational methods.

The FAIR guidelines are intentionally broad and somewhat vague; they are designed to provide communities with a flexible framework that can be further developed and adapted to specific needs. In practice, they focus primarily on metadata — covering aspects such as authorship, provenance, licensing, and descriptions of the dataset's contents.

Adopting FAIR principles can significantly improve the citability, visibility, and confirmability of datasets, making computational results more easily reproducible.

A brief note on interoperability: to the best of my knowledge, there is currently no standard knowledge representation language for mathematical data. While this means we don't yet have to worry about strict interoperability requirements, it is an area where the community should invest effort in the future.

If you think back to "everything is data" from earlier, a tricky question arises: can we state in general terms what metadata are sufficient to ensure reusability of data in mathematics?

The RDMP
Research Data Management Plan

- A living document, from the start of a project.
- Outlines how data will be managed throughout a research project.
- Incorporates the FAIR principles to ensure that data is handled in a way that maximizes its long-term value and usability.

Data Management Plan

No data management plan is necessary, since the research outlined in this proposal is in the realm of Mathematics and by nature theoretical. The PI will make his research freely and publicly accessible through articles, graphics and programs on the web page <INSERT WEBSITE>. The PI will also submit articles and papers to appropriate peer-reviewed journals for publication. Finally this work will be disseminated through academic research seminars and conferences.

www.math.harvard.edu/media/DataManagement.pdf
(link courtesy of Boege et al)

A claim I've often heard in and about the field of mathematics is that mathematicians rarely produce data, and that the data they do produce requires little to no management. I've also frequently come across statements like "you can't license mathematical objects" and the belief that if data is posted on someone's website, it is automatically in the public domain and freely usable.

A taste of an RDMP questionnaire

What we may end up dealing with if the current trends continue

- *Project metadata:* title, ID, grant reference, PI names, institution, ...
- *Expected questions:* RDMP author, data types and formats, how the data will be organized, secure storage and backup, documentation, volume of data, storage type, ...
- *Often disregarded:* license.
- *Possibly does not apply* to mathematics:
 - Ethics approval, legal issues, IP, culturally sensitive issues.
 - Data confidentiality and sensitivity, access restrictions (incl. cost)
 - Non-digital data questions
 - Data destruction

Documenting changes in the approach to collecting data can be informative
clear documentation of techniques is instrumental to reproducibility, also minimizes the
impact of onboarding new collaborators

Take-aways for data management

At this point in time and for people compiling datasets containing examples

- More and better metadata and documentation
though we are still waiting for a metadata standard
- Archiving and preservation: snapshots in a machine readable format on Zenodo or GitHub to ensure longevity
- Reproducibility for results of computation: record software info (version), attach code.
- An interesting problem up for community consideration: the meaning and provenance of mathematical data can require more complex mathematical data.
- A solution for recognition for research data beyond a journal publication is hard.

4. Where is the data?

We are still looking for the long tail

- An analysis of zbMATH references revealed a long tail of specific data
(Hulek, Müller, Schubotz and Teschke: Mathematical Research Data)
- *Call to action*
What are relevant metadata for collections of examples?

MATHBASES	
INDEX OF MATHEMATICAL DATABASES	
Info	Name
	(S-)Unit Groups of Orders ↗ Gabriele Nebe, Renaud Coulangeon, Oliver Braun, Sebastian Schönemann
	A Catalogue of Lattices ↗ Neil Sloane, Gabriele Nebe
	A Census of Edge-Transitive Tetravalent Graphs ↗ Primož Potočnik, Steve Wilson
	A Database for Number Fields ↗ Jürgen Klüners, Gunter Malle
	A Database of Continued Fractions of Polynomial Type ↗ Henri Cohen
	A Database of Galois polynomials ↗ Bill Allombert, Igor Schein
	A Database of Graphs in Combinatorica Format ↗ Sriram Pemmaraju, Steven Skiena
	Adjectives Project ↗ Jesse Voigl, David Holmes
	An Atlas of Abstract Regular Polytopes for Small Almost Simple Groups ↗ Laurence Vauthier, Dimitri Leemans
	An Atlas of Chiral Polytopes for Small Almost Simple Groups ↗ Dimitri Leemans, Michael Hartley, Isabel Hubard

MathBases began as an effort to index and showcase mathematical databases. David has already given you a tour of MathBases, so I won't dwell on that here — except to note that the strong focus on combinatorics is partly due to the relative approachability of combinatorial data, and partly a reflection of my own bias.

MathBases indexes datasets that contain examples of objects of interest to research mathematicians. When I was compiling its precursor MathDB, I applied a similar criterion, but I often struggled to decide whether or not a dataset should be included — even when I understood its contents.

As a call to action (building on David's list of ways you can contribute), I encourage you to think about what metadata are most relevant for datasets containing examples of mathematical objects.

The diversity of math databases

Four easy dimensions

small - large

Atlas of Small Chiral Polytopes (56) *Lists of finite lattices ($17 \cdot 10^9$)*

stored - generated on demand

House of Graphs - *The Small Groups Library in GAP - nauty*

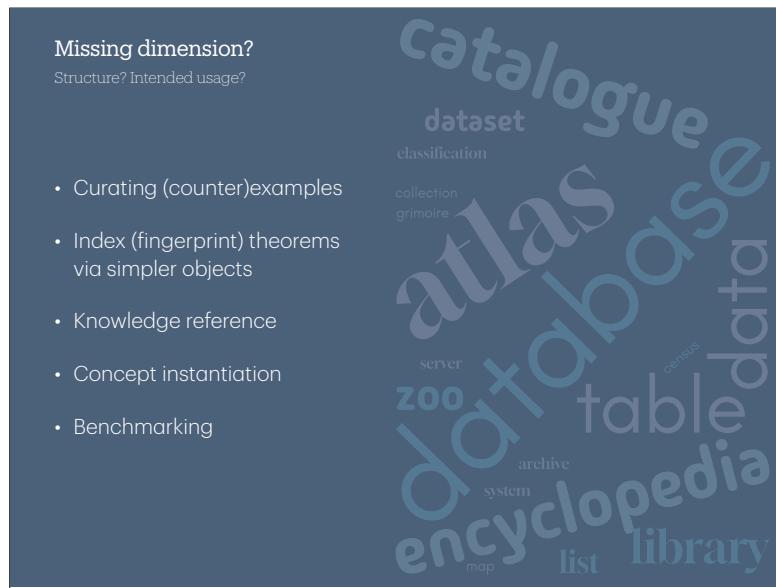
enumerated - curated/collected

The Foster census *House of Graphs*

easy to obtain values - a value corresponds to a paper

nauty *Parameters of Strongly Regular Graphs*

Room for more!



- Curation of examples: topological spaces, properties and theorems in π -base, graphs and invariants in the House of Graphs.
- **Index theorems: integer sequences (OEIS), Parameters of Strongly Regular Graphs.**
- Knowledge reference: definitions and properties of special functions in DLMF
- Instantiation: datasets in algebraic geometry (only one object, variety)
- Benchmarking: SuitSparse matrix dataset

A detour into the world of formalized mathematics: 1000+ theorems

Successor to Freek Wiedijk's 100 theorems

- 100 theorems (precursor): showcasing formalizations by keeping track of formalizations of the hundred greatest theorems (a *fixed* list)
- Indexing formalizations of a much longer (changing) list of theorems

Formalized theorems							
MSC	Name	Isabelle	HOL Light	Rocq	Lean	Metamath	Mizar
26	Abel's theorem		L				
12	Abel–Ruffini theorem		L				
68	Akra–Bazzi theorem		L				
11	Apery's theorem		X				
91	Arrow's impossibility theorem		X				
16	Artin–Wedderburn theorem		L				
46	Arzela–Ascoli theorem		L				

1000-plus.github.io

5. Trusting data

Beyond the standard checks

- How can we trust that a list of examples is complete (if applicable) and correct?
- Is the connection between theory and code sound?
- Are the results of computations correct?

No answers to the questions above, just an example:
Lean-HoG

$$\begin{aligned} A &= B \\ AA &= AB \\ AA - BB &= AB - BB \\ (A+B)(A-B) &= B(A-B) \\ A + B &= B \\ 2B &= B \\ 2 &= 1 \end{aligned}$$

Courtesy of Steven Clontz

Some standard options to increase the level of trust

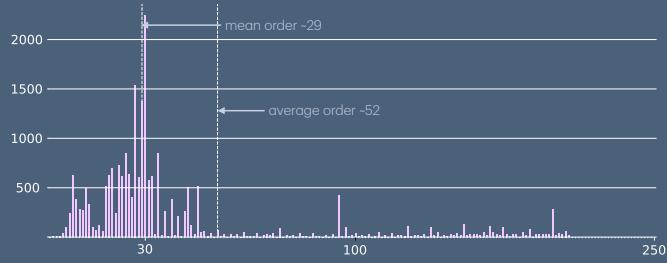
Techniques from engineering and mathematics

- Standard checks: format, type, consistency, uniqueness, ...
- *Testing*: software is run on a collection of test cases, the results are compared to reference results known to be true.
- *Redundancy*: several versions of software performing the same task are developed and executed independently, their results compared.
- Correctness of code or data is established by formal proofs.

The House of Graphs

The database of interesting graphs

- ~**32 000** graphs on up to 250 vertices,
- ~**50** properties, including computationally difficult ones, such as: genus, chromatic number, and Hamiltonicity.



The combination of graph sizes and properties means that we can't just compute whichever way we want.

Design options

based on quantity and complexity of objects and properties

- Prove the properties of each example by hand.
- Implement algorithm(s) in the proof assistant
(in the extreme case, implement a computer algebra system in a proof assistant)
- Encode as SAT, verify encoding to be correct, use a (trusted) solver, check the certificates provided by the solver.
- Use external software to compute properties *and their certificates*, use the proof assistant to check correctness.

1. For few objects and properties, simple.
2. Few properties, many objects, efficiently computable: can be difficult.
3. We used a combination of the last two.

Certificates (a.k.a. witnesses)

Instead of computing values, just check correctness

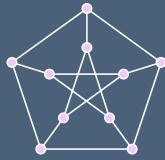
- Standard technique in computer science
- Check that 112909084933 is not a prime vs.
 $132241 \cdot 853813 = 112909084933$.

✓ Certificates *explain* why a property holds.

✗ Some properties do not have a certificate.

This works more broadly than you (might) think! (Used the idea for computation of election results).

The proof assistant can check the correctness of the certificate. While the connection with the property "not prime" follows directly from the definition here, this is not the case in general; a further proof that the property follows from the certificate can be necessary.



“The Petersen graph is a remarkable configuration that serves as a counterexample to many optimistic predictions about what might be true for graphs in general.”

Donald Knuth

Let's look at a random example of a graph from the House of Graphs.

The Petersen graph

has a Hamiltonian path, is not Hamiltonian



A foreshadowing of things to come.

Easy-ish: find a path in the graph that visits all vertices exactly once,
Harder: prove that we can't find such a cycle in the graph.

The Petersen graph is also the smallest vertex-transitive graph that is not a Cayley graph.

Lean-HoG

A Lean 4 library for finite simple graphs incorporating the House of Graphs

- Import graphs with efficient representations into Lean,
- together with values and certificates for:
the number of connected components, bipartiteness, traceability.
- A tactic to search the database and
- a tactic to close a goal by finding an example.
- Checking the number of connected components on (almost) all graphs takes ~16h.

Mathlib provides a basic, general-purpose formalization of simple graphs, but it was not suitable for our purposes. To address this, we implemented a small library for finite simple graphs, prioritizing efficiency over generality.

Early experiments showed that we could process a graph in time at most quadratic in the number of edges, and wherever possible, sub-quadratic in the number of vertices. Working naively with lists of vertices and edges — or with adjacency matrices — led almost immediately to quadratic (or worse) time complexity.

Some invariants, such as the number of edges, can be computed efficiently by the Lean kernel, provided an efficient graph representation. For other invariants — for example, testing bipartiteness via 2-coloring or detecting odd cycles — Lean can efficiently verify a certificate when supplied.

For the invariants (traceability), with certificates that only work in one direction, one strategy would be to complement them with heuristics wherever they work. For instance, detecting a disconnected graph is an easy way to rule out Hamiltonicity. Only when these simpler methods fail would we resort to SAT solving. However, we chose to take a more principled approach by using SAT for both directions.

Warning, implementation
details ahead.

Getting graphs into Lean

Mathlib: graphs represented with a symmetric, irreflexive adjacency relation

- Given a coloring c , check that adjacent vertices have different colors:
 $\forall ij : \text{Fin } n . \text{Adj } ij \rightarrow (ci \neq cj)$.
time complexity $\mathcal{O}(n^2)$, only $\mathcal{O}(|E|)$ when given a set of edges.
- Check whether a graph is regular:
 $\exists k : \mathbb{N} . \forall i : \text{Fin } n . |\{j : \text{Fin } n ; \text{Adj } ij\}| = k$
time complexity $\mathcal{O}(n^2)$, only $\mathcal{O}(n)$ when given a neighborhood map.

Lean-HoG graph representations

A Lean 4 library for finite simple graphs incorporating the House of Graphs

All properties require efficient

- membership checking, and
- checking that something holds for every element of a set, i.e. vertices, edges.

Lean-HoG:

- RBSet and RBMap for all sets and maps,
- graphs represented via sets of edges (and an auxiliary neighborhood map, checked to be equivalent).

Certificates:

could just to regular certificates (no SAT), with paths etc; for the other side use heuristics whenever they work (disconnected graph for Hamiltonicity), only resort to SAT when all else fails; we took the more principled approach with SAT.

Connected components

```
-- Verices u and v are connected if they are related by the equivalence
| | relation generated by the adjacency relation. -/
def Graph.connected {G : Graph} : G.vertex → G.vertex → Prop := EqvGen G.adjacent

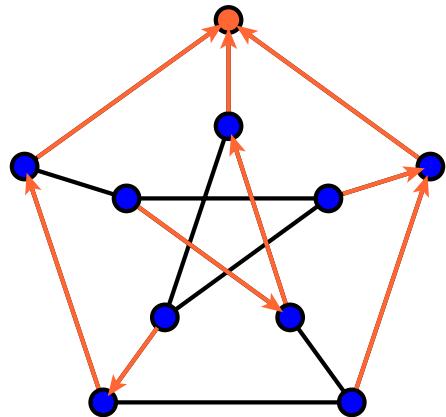
-- Connected components of a graph, as a structure -/
class ConnectedComponents (G : Graph) : Type :=
-- Number of connected components -/
val : Nat

-- The component of the given vertex -/
component : G.vertex → Fin val

-- Components are inhabited -/
componentInhabited : ∀ (i : Fin val), ∃ u, component u = i

-- The assignment of components coincides with connectedness -/
correct : ∀ u v, component u = component v ↔ G.connected u v
```

Certificate



Certificate

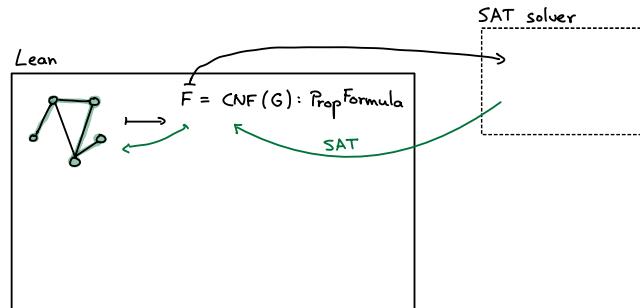
```
-- A certificate for connected components -/
class ConnectedComponentsCertificate (G : Graph) : Type :=
| -- Data
|   val : Nat
|   component : G.vertex → Fin val
|   root : Fin val → G.vertex
|   next : G.vertex → G.vertex
|   distToRoot : G.vertex → Nat
|
| -- Properties
|   componentEdge : G.edgeSet.all (fun e => component (G.fst e) = component (G.snd e)) =
|   rootCorrect : ∀ i, component (root i) = i
|   distRootZero : ∀ (i : Fin val), distToRoot (root i) = 0
|   distZeroRoot : ∀ (v : G.vertex), distToRoot v = 0 → v = root (component v)
|   nextRoot : ∀ i, next (root i) = root i
|   nextAdjacent : ∀ v, 0 < distToRoot v → G.adjacent v (next v)
|   distNext : ∀ v, 0 < distToRoot v → distToRoot (next v) < distToRoot v
|
-- From a components certificate we can derive the connected components -/
instance {G : Graph} [C : ConnectedComponentsCertificate G] : ConnectedComponents G :=
```

Load the certificate

```
--  
| JSON representation of connected components certificate.  
-/  
structure ConnectedComponentsData : Type where  
  val : Nat  
  component : Array (Nat × Nat)  
  root : Array (Nat × Nat)  
  next : Array (Nat × Nat)  
  distToRoot : Array (Nat × Nat)  
  
deriving Lean.FromJson  
  
-- Build a connected components certificate expression from the data. -/  
def buildCert (G : Q(Graph)) : ConnectedComponentsData → Q(ConnectedComponentsCertificate $G) :=
```

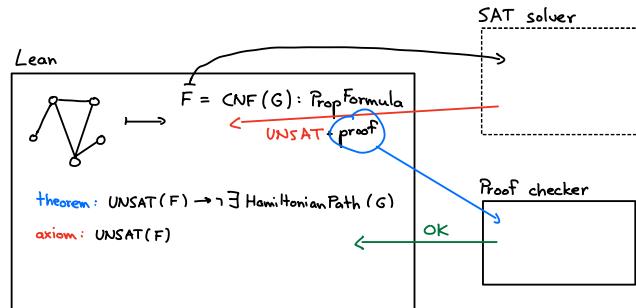
Hamiltonian Paths

NP-hard, in general no easy certificate for non-existence.
Use a SAT solver.



Hamiltonian Paths

NP-hard, in general no easy certificate for non-existence.
Use a SAT solver.



Take-aways for incorporating a database into a proof assistant

It probably won't be what you expected

- It depends on the database.
- Lean may be a sensible proof assistant to start with. If you do choose Lean, a lot depends on Mathlib.
- Checking a database may force you to consider efficiency and may make you feel like you are doing CS 50 years ago.
- Alternative to our approach: formal verification of algorithms.
- We implore database designers to consider certificates whenever possible.

We found it particularly advantageous to minimize the amount of computation performed directly by Lean, especially in situations involving meta-programming, where Lean metaprograms construct proofs for each value.

It would be possible to implement most of the properties of graphs in HoG. In some cases, however, we did not see a clear way out. For instance, computing the maximum or minimum eigenvalues of the adjacency matrix would require not only a standard format for algebraic numbers and a trusted, efficient computation engine for them, but also further considerations if we wanted to reason about extremality.

Thank you!

Recap of what we have been up to

- **MathBases**

Adam Towsley, Ben Spitz, David Roe, David Lowry-Duda, Benjamin Hutz, Edgar Costa, KB

- **Lean-HoG**

Jure Taslak, Gauvain Devillez, KB, Andrej Bauer

- **1000+ theorems**

Freek Wiedijk, Floris van Doorn, KB; editors for each system

- **MathDataHub**

Tom Wiesing, KB



Slides (hopefully) and references at katja.not.si, as promised



Andrej Bauer is looking
for someone to do a

Postdoc @Ljubljana

The mission:
*to create and curate large-scale
mathematical datasets and
dependency graphs extracted from
libraries of formalized mathematics.*

Start: **ASAP**
Talk to me for more info.

I feel like it is some kind of a rite of passage when you finally get to advertise a job. The PI on the project is Andrej Bauer (foundations of mathematics and logic, constructive and computable mathematics, homotopy type theory, mathematical foundations of programming languages, exact scientific computation, also, a very cool colleague). I will be working on the same project.

("can you lure a postdoc to Lj that would be more or less a copy of you")



Conference on Intelligent
Computer Mathematics

CICM 2026 @Ljubljana

Published proceedings,
peer-reviewed contributed papers,
database descriptions welcome

summer/September (TBD)
Hope to see you there!